

Offline Surgical QA with Decomposed Retrieval and Synthesis for Resource-Constrained Settings

Kiran Bhattacharyya

BHATTACHARYYAKIRAN12@GMAIL.COM *Intuitive Surgical, Inc., Atlanta, GA, USA*

Abstract

Digital access to critical medical knowledge in resource-limited settings is often hindered by a lack of internet connectivity and the computational demands of AI systems. This paper introduces the Surgical Information Assistant, a fully deployable, large language model (LLM) -driven multi-agent system designed to provide reliable surgical information in offline, resource-constrained environments. Our system is powered by a workflow that orchestrates question decomposition, information retrieval, grounded generation, and information synthesis to perform complex reasoning on consumer-grade hardware. Grounded in the Open Manual of Surgery for Resource-Limited Settings, we evaluated DeRetSyn on a new question-answer (QA) dataset of over 14,000 surgical question-answer pairs. We compare our system to other alternatives, perform ablation experiments on components of the agentic system, and interrogate sensitivity to retrieval parameters. The results show that our agentic orchestration enables a compact 3B Llama model to achieve 63% top-1 accuracy, significantly outperforming both a baseline GPT-4o (42.5%) and a larger 8B Llama model with conventional RAG (53%). We further test whether this performance enhancement from agentic orchestration for information retrieval generalizes to the PubMedQA dataset. Additionally, the entire system consumes <3.5 GB of RAM and generates responses within 8-15 seconds working on a consumer laptop. Our work serves as a practical blueprint for how agent-based systems can empower small, efficient models for medical domain information retrieval and synthesis, offering a tangible application of AI technology that could help advance health equity. We will release our dataset, code base, and prompts to foster further research in deployable and responsible clinical AI.

Keywords: clinical AI, information retrieval system, medical education, offline AI systems

Data and Code Availability The datasets used in this study are: 1) the PubMedQA dataset [Jin et al. \(2019\)](#) and 2) a dataset that we create and refer to as the *OpenManualOfSurgeryQA* dataset. The code and data is publicly available at [this Github repo](#). Please note that this repo is under active development and there may be breaking changes.

Institutional Review Board (IRB) This research does not require IRB approval since it is not “human subjects research” as it does not include activities that involve interaction with individuals or access to identifiable private information.

1. Introduction

Lack of access to modern technology, educational resource constraints, and limited opportunities for training have been identified as key challenges for surgical practice in resource-limited settings [Achanga et al. \(2025\)](#). Large language models (LLMs) offer a powerful way to democratize access to clinical knowledge with potential to improve surgical training and education. However, in many real-world clinical environments, especially in resource-constrained settings, the reliance of most LLMs on large-scale models, stable internet access, or high-performance hardware (e.g. GPUs) presents a significant barrier to adoption. Furthermore, even in clinical settings with advanced technology access, there are significant privacy concerns about sharing patient data with LLM service providers [Zhao and et al. \(2024\)](#). This work directly addresses these practical limitations by developing a tool designed to further equity in clinical AI by targeting environments where opaque reasoning, retrieval failures, lack of connectivity, and lack of patient privacy can hinder effective care [Ke et al. \(2025\)](#); [Ng et al. \(2025\)](#); [Shi et al. \(2025\)](#).

Agent-based systems that orchestrate multiple components to solve complex problems offer a promising paradigm for improving the reliability and addressing limitations of LLMs. A common approach is Retrieval-

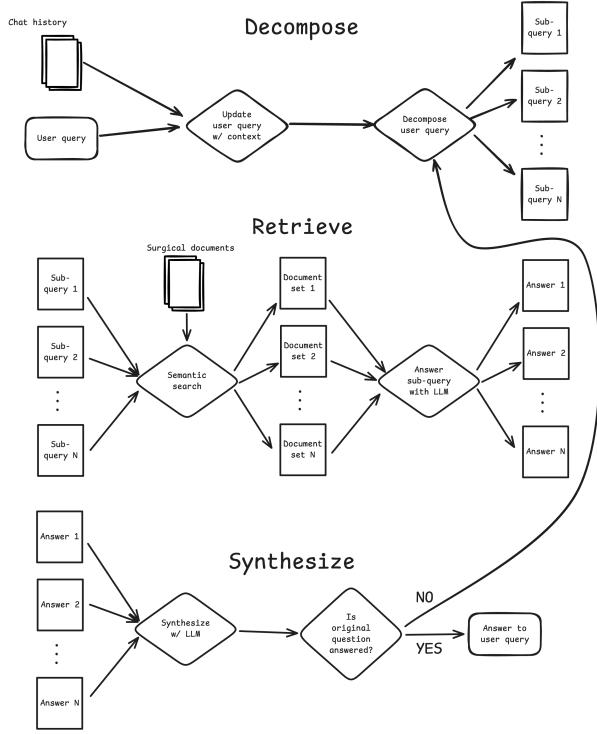


Figure 1: System design of the DeRetSyn pipeline.

Augmented Generation (RAG), which can be viewed as a simple two-stage workflow between a retriever and a generator agent to ground outputs in external knowledge [Vrdoljak et al. \(2025\)](#); [Bunnell et al. \(2025\)](#). More advanced systems now employ multi-agent, iterative workflows that decompose complex queries to emulate reasoning and improve interpretability [Zhang and et al. \(2024\)](#); [Xu et al. \(2024\)](#); [Liu et al. \(2025\)](#). However, these sophisticated agentic architectures often remain computationally intensive and reliant on cloud-based APIs, limiting their deployability where it is needed most.

To bridge this gap, we introduce the *Surgical Information Assistant*, an innovative and deployment-focused application designed for offline use. Built on the existing Open Manual of Surgery for Resource-Limited Settings [Vanderbilt University Medical Center \(2025\)](#), the system is designed to provide surgical information to user questions without an internet connection. It is powered by our multi-agent architecture, DeRetSyn (Decompose–Retrieve–Synthesize). This workflow orchestrates a collaboration between several specialized agents: 1) a DECOMPOSITION AGENT that breaks down complex surgical queries into answerable

sub-questions; 2) parallel RETRIEVAL AGENTS that fetch relevant text from a local, FAISS-indexed corpus; and 3) a SYNTHESIS AGENT that generates and iteratively refines answers under the guidance of a SUPERVISOR AGENT which verifies completion, using a compact Llama-3.2-3B model.

In an evaluation on a new dataset of over 14,500 QA pairs, our DeRetSyn-powered 3B model achieves 63.0% accuracy, significantly outperforming both a baseline GPT-4o (42.5%) and a larger 8B model using standard RAG (53.8%). Furthermore, the agentic orchestration shows a consistent advantage also on the PubMedQA dataset [Jin et al. \(2019\)](#). The system we describe here is deployable on consumer-grade hardware without internet connectivity.

While DeRetSyn is not formally affiliated with the Open Manual of Surgery, it builds on its content under the Creative Commons license. We are making our entire framework publicly available, including all code, data, and prompts, to support further research and development in accessible clinical AI. Our experiments also explicitly measure deployment relevant factors such as latency, memory, and offline usability, and we discuss dataset biases, generalizability, and practical limitations to inform safe clinical translation.

Our contributions are:

- A deployable, open-source agent-based system for offline surgical QA, demonstrating an application of LLM and AI agents for further equity in clinical AI.
- Empirical evidence that a multi-agent workflow can significantly boost the performance of small, efficient language models, making them viable for clinical QA in resource-constrained domains.
- A QA dataset to support further research and improvements in deployable agent-based systems for clinical AI.

2. Methods

2.1. Knowledge Base for a Low-Resource Context

The foundation of the Surgical Information Assistant is a curated, offline knowledge base derived from the *Open Manual of Surgery for Resource-Limited Settings (OMSRS)*. We selected this corpus for two primary reasons: 1) its content is explicitly tailored for clinicians practicing without access to specialized equipment or resources, directly aligning with our goal of promoting

health equity; and 2) its Creative Commons license permits redistribution, enabling the creation of a fully open-source and deployable system.

To prepare the corpus for retrieval, we parsed the source PDFs into text and segmented the content into chunks of 1000 characters with a 200-character overlap to maintain contextual continuity. These chunks were then converted into vector embeddings using the `all-MiniLM-L2-v2` model, chosen for its efficiency and strong performance in semantic search tasks. The resulting vectors were indexed into a local FAISS (Facebook AI Similarity Search) database, creating a fast and memory-efficient search index that can be deployed alongside the model on consumer-grade hardware and accessed via multiple threads. The total hard-disk space consumed by the index and the metadata was 60 MB.

2.2. The DeRetSyn Agentic Workflow

The core of the *Surgical Information Assistant* is DeRetSyn, a multi-agent workflow designed to emulate a structured reasoning process by breaking down complex problems into manageable steps (Figure 1). The system is coordinated by a lightweight, Python-based agentic controller that manages stepwise function calls for planning, retrieval, grounded generation, and synthesis. The workflow proceeds through the following specialized agents:

1. **Decomposition Agent:** An initial LLM-powered agent analyzes the user’s query to identify and extract underlying sub-questions. For example, a query like “How do I manage a post-operative wound infection?” might be decomposed into “What are the signs of a wound infection?”, “What are the common pathogens causing wound infections?”, and “What are the first-line treatments for wound infection in a low-resource setting?”.
2. **Parallel Retrieval Agents:** Each sub-question triggers an independent retrieval agent to perform a semantic search against the FAISS index. These agents run asynchronously to reduce overall response time, each retrieving the top- k most relevant text chunks for its assigned sub-question.
3. **Synthesis Agent:** For each sub-question and its retrieved context, a synthesis agent generates a concise, focused answer.
4. **Supervisor Agent:** A supervisor agent reviews the synthesized sub-answers to determine if the original query has been comprehensively addressed.

If gaps are identified, this agent generates new sub-questions and repeats the retrieval-synthesis loop, accumulating knowledge over a maximum of three iterations. This iterative refinement allows the system to correct its course and construct a more complete answer.

5. **Final Synthesis:** Once the supervisor agent confirms that the query has been fully answered, a final LLM call synthesizes all verified sub-answers into a single, coherent response for the user. All intermediate sub-questions and answers are also provided to the user for review and transparency. If the system cannot fully answer the query, it explicitly states this and lists uncovered sub-questions.

Fallback Mechanism for Out-of-Domain Queries. To handle queries that fall outside the scope of the primary surgical corpus, the system includes an optional, fallback mechanism. If the DeRetSyn workflow fails to find a satisfactory answer after three iterations, it can query a pre-indexed Wikipedia abstract dataset using ColBERTv2 as a secondary, “best effort” source [Khattab et al. \(2023\)](#). This component is disabled by default to ensure the core application remains fully functional offline without need for large data storage requirements (1.4 GB English abstracts 2020). For a complete set of prompts for all agents, please refer to the Appendix. Relevant portions of the Appendix are called out in-line throughout.

2.3. Deployment and Reproducibility

A core objective of this work is to produce a fully *deployable and reproducible* tool. The Surgical Information Assistant is implemented in Python and is designed to run entirely offline on consumer-grade hardware. The key components include:

- **Language Models:** We use the Llama-3.2-3B model as the primary LLM for all language generation, selected for its balance of performance in medical QA datasets and computational efficiency [Aaditya \(2024\)](#); [NVIDIA \(2024\)](#). For comparison, we also evaluate against GPT-4o and Llama-3.1-8B via API calls where noted. All inference was performed with a temperature of 0.
- **API and User Interface:** We provide an API-level and user-friendly web interface built with Streamlit, allowing for easy interaction with the system.

- **Open Source:** The entire codebase, including the FAISS index, prompts, and evaluation scripts, will be released publicly to encourage further research and adaptation.

For large-scale testing on QA datasets, the OpenAI API was used to call GPT-4o (March-August 2025) and the TogetherAI API was used to evaluate all other models (March-August 2025). DeRetSyn with the Llama-3.2-3B model was also benchmarked for runtime performance on consumer hardware with 8 CPU cores (Intel i7 10th Gen) and 16GB of RAM (no GPUs) running both Ubuntu 24.04 and Windows 10, as well as, a M2 Macbook Pro with 32GB RAM and a M3 Macbook Air with 16GB RAM both running Sequoia 15.5. For all testing, we use the 8-bit quantized version for Llama-3.1-8B and Llama-3.2-3B. We also test the 4-bit quantized version of Llama-3.2-3B for runtime and memory footprint benchmarking. We use applications like [LM Studio \(2024\)](#) or [Ollama \(2023\)](#) for local testing (based on operating system support) which wrap llama.cpp and use memory mapping to load model shards into RAM [Gerganov \(2023\)](#). We compare these results with community-reported benchmarks of Llama-3.2-3B on mobile devices.

3. Evaluation Framework

Our evaluation was designed to assess the Surgical Information Assistant across three key dimensions: 1) the quality and reliability of the novel QA dataset; 2) the accuracy and robustness of the DeRetSyn workflow compared to established baselines; and 3) the system’s practical performance on deployment-relevant metrics.

3.1. The OpenManualOfSurgeryQA Dataset

To facilitate reproducible research in low-resource clinical AI, we created and validated the *OpenManualOfSurgeryQA* dataset, a new dataset for surgical question-answering specific to resource-constrained settings.

Generation Process Starting with the OMSRS corpus, we divided the text into pseudo-randomly sized chunks (500 to 5000 characters) of varied context lengths and programmatically generated over 16,000 question-answer pairs using a Mistral Small 3.1 24B model to be answerable from these context chunks. The QA generation prompt consisted of several critical instructions. First, it mandates *groundedness*, requiring that answers be derived exclusively from the

provided text. Second, it requires *self-containment*, ensuring questions are fully understandable without the source passage. Third, it directs the model to create questions that are *research-able* having enough context that relevant information and passages can be retrieved to answer the question. Finally, the prompt guides the model to generate a diversity of question types (e.g., definition, causality, comparison, yes/no) and to synthesize answers concisely rather than merely extracting sentences verbatim. Lastly, the prompt mandates a step-by-step reasoning process and a structured output format. We also include a 1-shot example to help encourage instruction-following and format consistency.

Quality Curation To ensure the dataset’s quality and mitigate biases from automated generation, we implemented a multi-stage curation process. An initial manual review of 1,000 pairs by two individuals, including a practicing clinician, identified common failure modes (e.g., questions referencing figures or assuming external context). These patterns were used to create automated string-matching filters to remove low-quality pairs. Subsequently, each remaining QA pair was validated by an LLM against strict criteria for clarity, answerability, and faithfulness to the source text reducing the QA set further. Upon final manual review of another 1,000 pairs in the remaining set of 14,529, we did not discover any ambiguous or malformed QA pairs. This is the QA set we refer to as the *OpenManualOfSurgeryQA* dataset and use for further validation in this study.

3.2. Experimental Design

We conducted a series of experiments to answer three primary research questions:

1. **Efficiency vs. Scale:** Can a compact, offline model with an advanced retrieval workflow (DeRetSyn) outperform larger, state-of-the-art models that rely on parametric knowledge alone?
2. **Workflow Efficacy:** How does the multi-step DeRetSyn process compare against other RAG paradigms like vanilla retrieval and ReAct?
3. **Component Impact:** What are the specific contributions of the decomposition and fallback mechanisms to the system’s overall performance?

Comparative Models and Ablations To answer these questions, we benchmarked several configurations. On our *OpenManualOfSurgeryQA* dataset, we

compared DeRetSyn against baseline Llama-3.1-8B model, GPT-4o, vanilla RAG, and ReAct prompting Yao et al. (2023b) all with instructions for Chain-of-thought formatting. We also performed ablations by disabling the decomposition and fallback modules in DeRetSyn to isolate their impact.

To assess generalizability and mitigate the risk of evaluation circularity, we also evaluated DeRetSyn on the public *PubMedQA* dataset Jin et al. (2019) which contains “oracle” context for each question containing all relevant information. On this dataset, we compared our system against standard Chain-of-Thought (CoT) prompting Wei et al. (2022), both with and without provided context, to measure DeRetSyn’s ability to utilize evidence and reason effectively.

3.3. Metrics and Validation

Automated Accuracy Assessment We measured accuracy using a *LLM-as-judge* (Mistral Small 3.1 24B). The judge was prompted to evaluate if a generated answer was factually consistent, complete, and relevant compared to the ground-truth answer. To ensure reliability, the prompt included three few-shot examples of correct and incorrect evaluations.

Human Validation of the LLM-as-judge The validity of using an LLM-as-judge is critical for the credibility of our results. We conducted a blinded human evaluation on 100 randomly sampled model outputs (50 judged correct, 50 incorrect). A human evaluator, using the same criteria, showed a high degree of agreement with the LLM-as-judge, achieving a Cohen’s κ score of 0.90 (Table 1). This strong correlation provides confidence in our automated evaluation pipeline.

Table 1: Confusion matrix validating the LLM-as-judge against human evaluation for 100 QA pairs.

Human Label	LLM: Correct	LLM: Incorrect
Correct	48	3
Incorrect	2	47

We report 95% confidence intervals (CIs) for our accuracy metrics. Confidence intervals are computed either through the binomial approximation method for proportions of a population or through bootstrapping—we indicate which is used for all instances. To account for residual labeling uncertainty in the LLM-as-judge

(Cohen’s $\kappa = 0.90$ with human raters), we conservatively expand CIs by 10% to reflect the expected variance due to imperfect adjudication.

Deployment-Focused Benchmarking Beyond reporting accuracy on QA datasets, a primary goal of this work is practical utility in low-resource settings. Therefore, we also measured key deployment metrics for the DeRetSyn system. These metrics include:

- **End-to-End Latency:** The average time from query submission to final response.
- **Computational Cost:** The average number of tokens generated per query, which serves as a proxy for computational load.
- **Memory Footprint:** The peak RAM usage during inference to ensure the system can run on devices with limited memory.
- **Contextual Performance:** We report the performance of the system on those questions across all datasets that pertain to the most common operative procedures and pathologies encountered in resource-limited surgical settings, as described in Wong et al. (2014).

These practical benchmarks are essential for assessing the real-world viability of the Surgical Information Assistant.

4. Results

4.1. DeRetSyn Performance on Surgical QA

On the *OpenManualOfSurgeryQA* dataset, the DeRetSyn workflow with a 3B Llama model achieved a top-1 accuracy of 63.0%. As shown in Table 2, this result substantially outperforms both a baseline GPT-4o (42.5%) and a much larger 8B Llama model with vanilla RAG (52.5%). This ~10 percentage point improvement over standard RAG demonstrates the impact of the multi-step decomposition and refinement process.

Impact of Workflow Components Removing the initial decomposition step reduced accuracy by nearly 4 points to 59.2%, highlighting its importance in structuring the retrieval process. Note that the CIs of the full system compared with the ablated initial decomposition are *not* overlapping (Table 2). Disabling the fallback mechanism had a smaller impact (62.1% accuracy with overlapping CIs), which was expected since the evaluation questions were sourced from the primary corpus. The fallback to Wikipedia

Table 2: Accuracy on the OpenManualOfSurgeryQA dataset. DeRetSyn with a 3B model surpasses larger models and alternative RAG strategies. CIs are computed from the binomial approximation.

Model Configuration	Accuracy [95% CI] (%)
Llama-3.1-8B (base)	28.7 [26.9, 30.4]
GPT-4o (base)	42.5 [40.7, 44.3]
Llama-3.1-8B + Vanilla RAG	52.5 [50.7, 54.3]
Llama-3.2-3B + Vanilla RAG	53.8 [52.0, 55.6]
Llama-3.2-3B + ReACT RAG	53.5 [51.7, 55.3]
Llama-3.2-3B + DeRetSyn (Ablated: Decomposition)	59.2 [57.4, 61.0]
Llama-3.2-3B + DeRetSyn (Ablated: Fallback)	62.1 [60.3, 63.9]
Llama-3.2-3B + DeRetSyn (Full System)	63.0 [61.2, 64.8]

was triggered in only 5% of cases for the full system, confirming its role as a supplementary, rather than essential, component for in-domain queries.

Performance Stability A bootstrap analysis of accuracy from 200 subsets containing 30 samples each (Figure 2) shows that DeRetSyn has the highest mean accuracy but also exhibits stable performance with a standard deviation similar to other configurations. From bootstrapping, the performance of the full DeRetSyn system was significantly different from all configurations ($p < 0.001$, Unpaired t-tests with Bonferroni correction) except for DeRetSyn with ablated fallback.

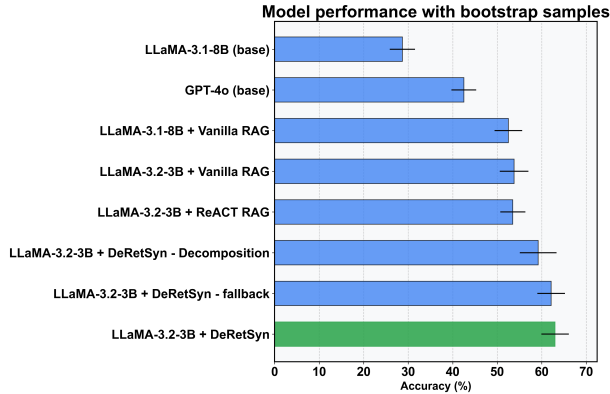


Figure 2: Bootstrapped accuracy comparison. DeRetSyn demonstrates higher mean accuracy and comparable variance to other configurations.

Computational cost We also computed the average number of tokens generated per response (Figure

3). Interestingly, Llama-3.2-3B generally used fewer tokens than the other LLMs tested in this study under similar configurations. However, the DeRetSyn system generates a larger number of total tokens than other RAG configurations due to the multiple generations required for each final response. Notably, the rate of increase in token generation is still less than expected when compared to the Llama-3.1-8B and GPT-4o base models under CoT prompting. While total generated tokens across all LLM calls for the full DeRetSyn system averaged ~ 600 , the final synthesized answers shown to users averaged ~ 100 tokens.

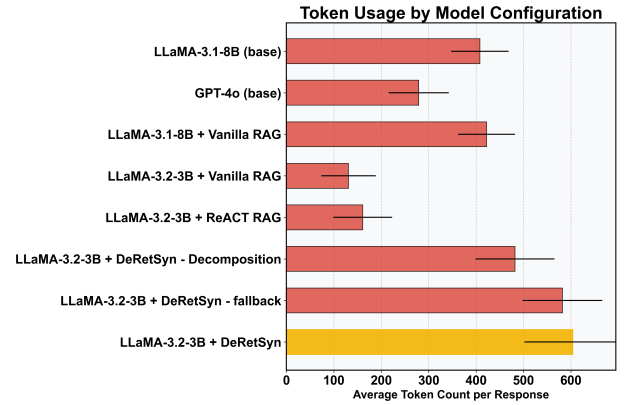


Figure 3: Total number of tokens generated from each configuration evaluated. Error bars represent the standard deviation in the dataset. Total tokens generated are not equal to the number of tokens in the final answer presented to the user due to CoT prompting and intermediate tool-usage steps.

4.2. Generalizability on the PubMedQA Dataset

To validate that DeRetSyn’s effectiveness was not limited to our own dataset, we evaluated it on the public PubMedQA dataset. As shown in Table 3, DeRetSyn with a provided context achieves **70.1%** accuracy, outperforming both a similarly prompted GPT-4o (64.1%) and standard CoT prompting (63.5%). This demonstrates DeRetSyn’s superior ability to utilize provided evidence for reasoning.

Notably, even when forced to rely only on its internal reasoning structure (*no context, no retrieval*), DeRetSyn (52.1%) still outperforms a baseline CoT

model (50.0%), suggesting the decomposition process itself aids reasoning.

4.3. Practical Performance for Real-World Deployment

Beyond theoretical accuracy, the practical viability of the Surgical Information Assistant depends on its performance on consumer-grade hardware and its ability to support the most common operative procedures and pathologies encountered in resource-limited settings.

Memory Footprint and Feasibility on Consumer Hardware Memory consumption is an important factor for hardware deployment in low-resource settings. We benchmarked the Surgical Information Assistant on a range of consumer-grade hardware to measure its peak RAM usage during inference. As shown in Table 4, the peak memory usage for the Llama-3.2-3B model was consistently under 3.5 GB across all platforms. This is lower than expected due to memory mapping that is handled by llama.cpp by only loading necessary parts of the model into memory. This low memory footprint confirms that the entire system can run on standard laptops and devices with as little as 8 GB of total RAM. The results also show a modest latency improvement and lower memory consumption with 4-bit quantization.

Performance on Relevant Procedures and Pathologies We automatically classified all questions in the *PubMedQA* and *OpenManualOfSurgeryQA* datasets with zero-shot classification using the `bart-large-mnli` model to determine which ones were relevant to the most common operative procedures and pathologies in resource-limited settings as reported in Wong et al. (2014). These included ‘Cesarean delivery’, ‘herniorrhaphy’, ‘wound debridement and care’, ‘percutaneous infections and abscesses’, and ‘circumcision’. We also included an ‘other’ category for those questions that did not pertain to any of the above. We found a total of 1372 questions across all relevant categories. The DeRetSyn system had a higher accuracy for this subset of questions (75.1%, 96% CI: [72.7, 77.3]) than for the entire *OpenManualOfSurgeryQA* data.

4.4. Analysis of Failure Modes

A qualitative review of incorrect answers revealed three primary failure modes:

- **Sub-question Misalignment:** Occurs when the decomposition agent misinterprets the intention of

a query, leading to sub-questions that are too broad or tangential.

- **Synthesis Drift:** In cases of sparse or conflicting retrieved evidence from nuanced questions (e.g., PubMedQA ‘Maybe’ items), the synthesis agent occasionally blended factual information with plausible-sounding fabrications. These errors are usually visible in intermediate steps, supporting safe human oversight.
- **Generic Fallback:** When the fallback to Wikipedia was triggered, the retrieved information sometimes lacked the necessary clinical specificity. These findings underscore the need for ongoing work in prompt engineering for clinical nuance, improving synthesis grounding, and curating more domain-specific fallback resources.

5. Discussion

Our findings demonstrate that with the an appropriate architecture, compact, locally-runable models can outperform cloud-based systems in specialized domains. This work provides a practical blueprint for developing LLM-based tools that are transparent and deployable. While the system described clearly has limitations, we believe it presents an architecture that aligns with the technical infrastructure requirements of a resource-constrained clinical setting and improves over baselines and much larger models without extensive fine-tuning. Each component of this architecture can be further improved while staying aligned with the core requirements for a resource-limited setting.

5.1. Bridging the Performance Gap with Efficient Orchestration

A central result of our study is that knowledge access and intelligent workflow orchestration are more critical than raw model scale for specialized tasks. The Llama-3.2-3B model, when augmented with the DeRetSyn workflow, surpassed the much larger GPT-4o and a vanilla RAG implementation with an 8B model for the datasets we tested. This confirms that the performance gap between large and small models can be effectively closed by shifting focus from scaling parametric knowledge to optimizing the retrieval and synthesis process Singhal and et al. (2023); Wang and et al. (2024).

The DeRetSyn method’s 10.5 percentage point accuracy gain over standard RAG highlights the value of

Table 3: Performance on the PubMedQA dataset. DeRetSyn demonstrates superior reasoning and evidence utilization compared to baselines. CIs are computed from the binomial approximation.

Method	Overall [95% CI] (%)	‘Yes’ (%)	‘No’ (%)	‘Maybe’ (%)
1. LLM-only w/o CoT and w/o context	48.50 [45.4, 51.7]	46.92	63.31	10.91
2. LLM-only w/ CoT and w/o context	50.00 [46.9, 53.2]	66.12	35.80	12.73
3. LLM-only w/ CoT and fixed PubMedQA context	63.50 [60.4, 66.4]	71.56	66.57	13.64
4. DeRetSyn (no context, no fallback, no retrieval)	52.10 [49.0, 55.2]	67.39	39.94	12.73
5. DeRetSyn w/ fixed PubMedQA context (no fallback, no retrieval)	70.10 [67.2, 72.9]	79.89	71.89	15.45
6. DeRetSyn w/ OMSRS retrieval and fallback (no PubMedQA context)	61.50 [58.4, 64.5]	68.84	65.09	13.63
7. <u>GPT-4o</u> CoT w/ fixed PubMedQA context	64.10 [61.0, 67.1]	60.87	68.05	68.18

Table 4: Performance and memory benchmarks on consumer hardware of the DeRetSyn-powered Llama-3.2-3B model.

Hardware Platform	Quantization	Avg. Latency (s)	Peak RAM (GB)
Windows/Ubuntu (8 CPU, 16GB RAM)	8-bit	14.2 - 15.6	3.1 - 3.2
	4-bit	13.5 - 13.8	2.5
M2 Macbook Pro (32GB RAM)	8-bit	8.2	2.6
	4-bit	7.7	2.2
M3 Macbook Pro (16GB RAM)	8-bit	8.1	2.5
	4-bit	7.7	2.1

its structured, multi-step approach. By decomposing complex queries, the system emulates a more deliberate reasoning process, which, as our PubMedQA results show, improves both the utilization of provided context and the model’s intrinsic reasoning capabilities.

5.2. Implications for Health Equity and Responsible AI

The practical impact of this work is rooted in its potential to advance health equity. The Surgical Information Assistant is designed to be fully operational on consumer-grade hardware without an internet connection, directly addressing the digital divide that excludes many low-resource clinical settings from the benefits of modern AI. A preliminary analysis shows that the system handles queries about common procedures and pathologies for resource-limited settings with a higher accuracy than most queries overall. However, this requires further investigation within more clinically relevant settings. Critically, the current performance of the system (~60%–70% accuracy) is likely not adequate for immediate use by clinicians in the field. We discuss how this can be improved further.

Trust and Interpretability in a Clinical Workflow Beyond accessibility, the DeRetSyn architecture promotes trust and transparency. DeRetSyn’s

modularity makes the reasoning process explicit; by surfacing the intermediate sub-questions, retrieved evidence, and synthesized sub-answers, it allows a clinician to quickly verify the basis of a recommendation. Moreover, our system can localize patient data since it can run on-device enabling data security. This aligns with principles of interpretable AI [Wang et al. \(2022\)](#); [Yao et al. \(2023a\)](#) and is essential for safe and responsible deployment in medicine [Zhao and et al. \(2024\)](#). In future studies, we plan to leverage this transparent reasoning to study how this tool can improve clinical judgment.

The Accuracy vs. Efficiency Tradeoff DeRetSyn’s accuracy comes at the cost of increased token generation (Figure 3). However, our benchmarks on laptops confirm that this latency remains well within practical limits. On mobile devices, the Llama-3.2-3B model with a llama cpp deployment generates 17–23 tokens/second [Arm Newsroom \(2024\)](#); [Qualcomm \(2024\)](#); [Hotellnx \(2024\)](#); [ctrl-brk \(2024\)](#). Assuming an average of 600 tokens for the final response (Figure 3), latency per query is estimated at ~30 seconds. Devices capable of running multiple model instances in parallel may reduce this time through concurrent execution. Moreover, user interactions with the intermediate steps of the system may be

important for transparency and auditability, making the total latency not the sole focus of usability.

5.3. Limitations and Future Directions

While our results are promising, we identify many areas for future work:

Dataset and Evaluation: The primary evaluation relied on a dataset generated from the source corpus. Although a strong performance on PubMedQA mitigates concerns of circularity, future work should validate the system against a dataset created from real-world clinical queries. Furthermore, future QA evaluations should incorporate more complex, multi-hop reasoning scenarios.

Prospective Clinical Evaluation: The ultimate test of this system is its performance in a real clinical setting. We are actively planning a prospective study to evaluate the Surgical Information Assistant with clinicians in a simulated low-resource environment. This will be crucial for assessing usability, measuring its impact on decision-making, and identifying unforeseen failure modes or risks.

Human-in-the-Loop Refinement: The current DeRetSyn workflow is fully automated. A promising future direction is to incorporate human-in-the-loop capabilities, allowing a clinician to correct a flawed sub-question or provide additional context during the refinement process. This could further enhance accuracy and usability.

Multi-language support: While DeRetSyn is language-agnostic, the OMSRS corpus is in English which constrains the current language of the reference documentation. We tested a 100-question subset of the *OpenManualOfSurgeryQA* translated to Spanish, showing comparable accuracy after back-translation of the final answer with confidence intervals overlapping with of the full system (-1.5% drift of mean accuracy). Future work will investigate this further and integrate multilingual front-ends and bilingual retrieval indexes for broader accessibility.

Knowledge and confidence: While DeRetSyn achieves strong results relative to the size of the LLM, its reasoning and retrieval accuracy are ultimately bounded by the pretrained model’s latent knowledge and its ability to assess uncertainty. Hypothetically, both could be addressed through fine-tuning and confidence modeling efforts. Specifically, curating “good” DeRetSyn traces from multiple medical QA datasets could be used to finetune smaller LLMs for improved domain adaptation and reasoning stability. In parallel,

developing internal confidence estimators at each sub-question and synthesis stage to quantify uncertainty and surface low-confidence outputs to users could improve transparency and performance. Together, these efforts could strengthen both the factual grounding and epistemic transparency of the system.

Potential harms: Potential harms include overconfident or outdated recommendations, especially under sparse evidence. The system can partially mitigate this through explicit deferrals (‘insufficient evidence’) and the visible citations per sub-question. Future work includes uncertainty-aware confidence scores and clinician-in-the-loop verification during the initial decomposition.

6. Conclusion

We show how the Surgical Information Assistant addresses resource-constraints and transparency by ensuring offline access to relevant medical information through explicit iterative reasoning and synthesis with reasonable latency.

Acknowledgments

We acknowledge Dr. Sankha Banerjee for their support in helping to review question-answer pairs and system responses.

References

- K Aaditya. Performance Comparison: Llama-3.2 vs. Llama-3.1 LLMs and Smaller Models (3B, 1B) in Medical and Healthcare AI Domains. Hugging Face Blog, Sep 2024. URL <https://huggingface.co/blog/aaditya/llama3-in-medical-domain>.
- BillSmith Anyinkeng Achanga, Christian Wabene Bisimwa, Victor Oluwafemi Femi-Lawal, Nnoko Sona Akwo, and Tohson Falake Toh. Surgical practice in resource-limited settings: Perspectives of medical students and early career doctors: A narrative review. *Health Science Reports*, 8(1):e70352, 2025.
- Arm Newsroom. Ai inference everywhere with new llama llms on arm. <https://newsroom.arm.com/news/ai-inference-everywhere-with-new-llama-llms>, 2024. URL <https://newsroom.arm.com/news/>

- ai-inference-everywhere-with-new-llama-llms. LM Studio. Lm studio documentation, 2024. URL <https://lmstudio.ai/docs/app>. Accessed: 2024-08-22.
- David J Bunnell, Mary J Bondy, Lucy M Fromtling, Emilie Ludeman, and Krishnaj Gourab. Bridging ai and healthcare: A scoping review of retrieval-augmented generation—ethics, bias, transparency, improvements, and applications. *medRxiv*, pages 2025–04, 2025.
- ctrl-brk. Phone LLM’s benchmarks? Reddit, Nov 2024. URL https://www.reddit.com/r/LocalLLaMA/comments/1glx6a5/phone_llms_benchmarks/.
- Georgi Gerganov. llama.cpp. <https://github.com/ggml-org/llama.cpp>, 2023.
- Hotellnx. Run Llama 3.2 3B on Phone - on iOS & Android. Reddit, Oct 2024. URL https://www.reddit.com/r/LocalLLaMA/comments/1fppt99/run_llama_32_3b_on_phone_on_ios_android/.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. PubMedQA: A Dataset for Biomedical Research Question Answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2201–2211, 2019.
- Yu He Ke, Liyuan Jin, Kabilan Elangovan, Hairil Rizal Abdullah, Nan Liu, Alex Tiong Heng Sia, Chai Rick Soh, Joshua Yi Min Tung, Jasmine Chiat Ling Ong, Chang-Fu Kuo, et al. Retrieval augmented generation for 10 large language models and its generalizability in assessing medical fitness. *npj Digital Medicine*, 8(1):187, 2025.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, et al. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*, 2023.
- Siru Liu, Allison B McCoy, and Adam Wright. Improving large language model applications in biomedicine with retrieval-augmented generation: a systematic review, meta-analysis, and clinical development guidelines. *Journal of the American Medical Informatics Association*, page ocaf008, 2025.
- Karen Ka Yan Ng, Izuki Matsuba, and Peter Chengming Zhang. Rag in health care: a novel framework for improving communication and decision-making by addressing llm limitations. *NEJM AI*, 2(1):AIra2400380, 2025.
- NVIDIA. Llama-3.2-3b-instruct Model by Meta - NVIDIA NIM APIs. <https://build.nvidia.com/meta/llama-3.2-3b-instruct/modelcard>, 2024. Accessed: August 4, 2025.
- Ollama. Ollama. <https://github.com/ollama/ollama>, 2023.
- Qualcomm. qualcomm/Llama-v3.2-3B-Instruct. Hugging Face, Sep 2024. URL <https://huggingface.co/qualcomm/Llama-v3.2-3B-Instruct>.
- Yucheng Shi, Shaochen Xu, Tianze Yang, Zhengliang Liu, Tianming Liu, Xiang Li, and Ninghao Liu. Mkrag: Medical knowledge retrieval augmented generation for medical question answering. In *AMIA Annual Symposium Proceedings*, volume 2024, page 1011, 2025.
- K Singhal and et al. Retrieval-augmented language models for clinical medicine. *Nature*, 620:282–289, 2023.
- Vanderbilt University Medical Center. About the global surgical atlas. <https://www.vumc.org/global-surgical-atlas/about>, 2025. Accessed: 2025-04-02.
- Josip Vrdoljak, Zvonimir Boban, Marino Vilović, Marko Kumrić, and Joško Božić. A review of large language models in medical education, clinical decision support, and healthcare administration. In *Healthcare*, volume 13, page 603. MDPI, 2025.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Z Wang and et al. Medgraphrag: Safe medical llm via graph-based retrieval-augmented generation. *arXiv preprint arXiv:2408.04187*, 2024.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

Evan G Wong, Miguel Trelles, Lynette Dominguez, Shailvi Gupta, Gilbert Burnham, and Adam L Kushner. Surgical skills needed for humanitarian missions in resource-limited settings: common operative procedures performed at medecins sans frontieres facilities. *Surgery*, 156(3):642–649, 2014.

Kehan Xu, Kun Zhang, Jingyuan Li, Wei Huang, and Yuanzhuo Wang. Crp-rag: A retrieval-augmented generation framework for supporting complex logical reasoning and knowledge planning. *Electronics*, 14(1):47, 2024.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023. *URL* <https://arxiv.org/abs/2305.10601>, 3:1, 2023a.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023b.

Hao Zhang and et al. Discuss-rag: Enhancing retrieval-augmented generation via agent-led multi-turn reasoning for medical qa. *arXiv preprint arXiv:2504.21252*, 2024.

X Zhao and et al. Aipatient: Simulating patients with ehers and agentic workflow. *arXiv preprint arXiv:2409.18924*, 2024.

Appendix A. Instruction prompt for QA generation from text fragments

You are an expert medical curriculum developer. Given a medical or medically-related text passage, you must generate question-answer pairs. The questions MUST be

- self-contained: answerable by the information in the passage but understandable on their own without needing direct reference to the passage (ie. "what is this passage about?" is NOT a valid question).
- grounded: directly related to the information in the passage
- research-able: have enough information to be able to retrieve the relevant passages or information in the future to help answer the question.

Generate a set of question-answer types including yes/no answer, short phrases, and longer sentences.

Think step-by-step and reason through the content of the passage to hypothesize potential question-answer pairs and then respond. Think of at least 8 question-answer pairs but more, if possible.

Respond in this format:

```
<think> You reasoning here... </think>
<qa-pair> [question here...] | [answer here...] </qa-pair>
<qa-pair> [question here...] | [answer here...] </qa-pair>
```

Here is an example:

PASSAGE:

The goal of transvesical prostatectomy is to remove the hypertrophic transitional zone of the prostate to reduce urinary obstruction. It is not a surgical procedure for malignant lesions. Appropriate history (including functional

status inquiry) and physical examination are done prior to any investigation. Physical exam always includes digital rectal exam. Enlarged lateral lobes can be easily felt and the rectal mucosa should move freely over the prostate. Any nodularity or fixation of the rectal mucosa should alert you to the possibility of prostatic cancer and core needle biopsy should be done before proceeding with the admission.

RESPONSE:
<think>

Okay, so I need to generate question-answer pairs based on the given medical passage. The questions should be answerable from the passage and understandable on their own. Let me read through the passage carefully to understand the key points.

The passage is about transvesical prostatectomy. It mentions that the goal is to remove the hypertrophic transitional zone of the prostate to reduce urinary obstruction. It's not for malignant lesions. Before any investigation, they do a history and physical exam, including a digital rectal exam. Enlarged lateral lobes can be felt, and the rectal mucosa should move freely. If there's nodularity or fixation, it might indicate cancer, so a biopsy is needed before admission.

Alright, so the main topics here are the purpose of the procedure, what it's not used for, the pre-procedure assessments, the physical exam findings, and when to suspect cancer.

I should think of questions that cover these areas. Let's start with the purpose. A question like, "What is the primary goal of transvesical prostatectomy?" The answer would be about removing the hypertrophic zone to relieve obstruction.

Next, since it's not for malignant lesions, a question could be, "Is transvesical prostatectomy used to treat malignant lesions of the prostate?" The answer would be no.

Then, the pre-procedure steps. Maybe, "What assessments are typically performed before transvesical prostatectomy?" The answer would include history, physical exam, and digital rectal exam.

For the physical exam findings, a question could be, "What is a normal finding during the digital rectal exam in a patient undergoing transvesical prostatectomy?" The answer would mention enlarged lateral lobes and freely moving rectal mucosa.

When to suspect cancer: "What signs during a digital rectal exam may indicate the need for a core needle biopsy?" The answer would be nodularity or fixation of the rectal mucosa.

I can also think about the implications of the findings. For example, "What should be done if nodularity or fixation is detected during the digital rectal exam?" The answer would be to perform a biopsy before proceeding.

Another question could be about the importance of the digital rectal exam in this context. "Why is a digital rectal exam important before transvesical prostatectomy?" The answer would be to assess for signs of cancer.

I should make sure each question is clear and can be answered without directly referencing the passage, just the information in it. Let me see if I can come up with more. Maybe about the specific zones of the prostate involved. "Which part of the prostate is targeted in transvesical prostatectomy?" The answer is the transitional zone.

Also, the procedure's focus: "What is the focus of transvesical prostatectomy in treating urinary obstruction?" The

974	answer is removing the hypertrophic	a core needle biopsy? Nodularity or	1025
975	transitional zone.	fixation of the rectal mucosa may	1026
976	I think that's a good set. I have eight	indicate the possibility of prostatic	1027
977	pairs, but maybe I can find a couple	cancer. </qa-pair>	1028
978	more. Perhaps about the consequences	<qa-pair> What should be done if	1029
979	of not doing the biopsy. "What is the	nodularity or fixation is detected	1030
980	consequence of not performing a core	during the digital rectal exam? A	1031
981	needle biopsy when suspicious findings	core needle biopsy should be performed	1032
982	are present?" The answer would be	before proceeding with the admission.	1033
983	proceeding without knowing if cancer	</qa-pair>	1034
984	is present, which could be risky.	<qa-pair> Why is a digital rectal exam	1035
985	Wait, the passage doesn't explicitly state	important before transvesical	1036
986	the consequences, so maybe that's	prostatectomy? It is important to	1037
987	beyond the given info. I should stick	assess for signs of prostatic cancer,	1038
988	to what's in the passage.	such as nodularity or fixation of the	1039
989	Another angle: the patient's history.	rectal mucosa. </qa-pair>	1040
990	"What is included in the patient's	<qa-pair> Which part of the prostate is	1041
991	history before transvesical	targeted in transvesical	1042
992	prostatectomy?" The answer would be	prostatectomy? The transitional zone	1043
993	functional status inquiry.	of the prostate is targeted. </qa-pair>	1044
994	I think that's a solid list. I'll make	<qa-pair> What is the focus of	1045
995	sure each question is clear and each	transvesical prostatectomy in treating	1046
996	answer is accurate based on the	urinary obstruction? The focus is on	1047
997	passage.	removing the hypertrophic transitional	1048
998	</think>	zone of the prostate. </qa-pair>	1049
999	<qa-pair> What is the primary goal of	<qa-pair> What is included in the	1050
1000	transvesical prostatectomy? The	patient's history before transvesical	1051
1001	primary goal is to remove the	prostatectomy? The patient's history	1052
1002	hypertrophic transitional zone of the	includes a functional status inquiry.	1053
1003	prostate to reduce urinary	</qa-pair>	1054
1004	obstruction. </qa-pair>		1055
1005	<qa-pair> Is transvesical prostatectomy		1056
1006	used to treat malignant lesions of the	Here is the passage for which you need to	1057
1007	prostate? No, it is not a surgical	generate question-answer pairs	1058
1008	procedure for malignant lesions.		1059
1009	</qa-pair>	PASSAGE:	1060
1010	<qa-pair> What assessments are typically	{text_chunk}	1061
1011	performed before transvesical		1062
1012	prostatectomy? Appropriate history,	RESPONSE:	1063
1013	including functional status inquiry,		
1014	and physical examination are done		
1015	prior to any investigation. </qa-pair>		
1016	<qa-pair> What is a normal finding during	Appendix B. System prompt for	1064
1017	the digital rectal exam in a patient	LLM-as-judge	1065
1018	undergoing transvesical prostatectomy?	### ROLE AND OBJECTIVE ###	1066
1019	Enlarged lateral lobes can be easily	You are an impartial and meticulous AI	1067
1020	felt, and the rectal mucosa should	evaluator. Your objective is to	1068
1021	move freely over the prostate.	determine if a "Generated Answer"	1069
1022	</qa-pair>	correctly and completely answers a	1070
1023	<qa-pair> What signs during a digital	given "Query", using the "Ground Truth	1071
1024	rectal exam may indicate the need for	Answer" as the definitive source of	1072

1073	correctness. Your evaluation must	Answer" based on the **Factual	1121
1074	result in a binary decision: "Correct"	Consistency** and **Completeness**	1122
1075	or "Incorrect".	criteria.	1123
1076		4. **Final Judgment:** Based on your	1124
1077	### EVALUATION CRITERIA ###	comparison, make a final binary	1125
1078	You must adhere to the following strict	judgment.	1126
1079	criteria:	5. **Output Format:** Provide your	1127
1080		response in the following format	1128
1081	1. **Factual Consistency:** The "Generated	<think> Your reasoning here... </think>	1129
1082	Answer" must be factually consistent	<answer> correct OR incorrect </answer>	1130
1083	with the "Ground Truth Answer". It		1131
1084	must not contain any information that	Here are some examples	1132
1085	contradicts the ground truth.		1133
1086	2. **Completeness:** The "Generated	**Example 1**	1134
1087	Answer" must address all parts of the		1135
1088	"Query". It is considered "Incorrect"	QUERY	1136
1089	if it omits critical information that	What is the escape velocity from the	1137
1090	is present in the "Ground Truth	surface of Earth?	1138
1091	Answer" and is necessary for a full		1139
1092	response.	GROUND TRUTH ANSWER	1140
1093	3. **Relevance:** The "Generated Answer"	The escape velocity from Earth's surface	1141
1094	must directly answer the user's	is approximately 11.2 kilometers per	1142
1095	"Query". An answer that is factually	second (km/s), which is about 40,320	1143
1096	correct but irrelevant to the question	kilometers per hour (km/h) or 25,000	1144
1097	is "Incorrect".	miles per hour (mph).	1145
1098			1146
1099	**Important Note:** Differences in	GENERATED ANSWER	1147
1100	phrasing, verbosity, or style between	To escape Earth's gravity from the	1148
1101	the "Generated Answer" and the "Ground	surface, an object needs to reach a	1149
1102	Truth Answer" are acceptable as long	speed of about 11.2 km/s.	1150
1103	as the core semantic meaning is the		1151
1104	same and the criteria above are met.	YOUR RESPONSE	1152
		<think> Step 1: The query asks for Earth's	1153
1105	B.1. Instruction prompt prefix for	escape velocity. Step 2: The ground	1154
1106	LLM-as-judge	truth states it is 11.2 km/s. Step 3:	1155
		The generated answer also states it is	1156
1107	### INSTRUCTIONS ###	11.2 km/s. Step 4: The generated	1157
1108	Follow these steps to generate your	answer is factually consistent and	1158
1109	evaluation:	directly answers the query. While less	1159
1110		detailed than the ground truth, it	1160
1111	1. **Analyze the Query:** Carefully read	does not omit the critical piece of	1161
1112	the user's "Query" to understand the	information and contains no	1162
1113	user's full intent and all a-priori	contradictions. Therefore, it is	1163
1114	information.	correct. </think>	1164
1115	2. **Study the Ground Truth:** Read the	<answer> correct </answer>	1165
1116	"Ground Truth Answer" to understand		1166
1117	the complete and factually correct	**Example 2**	1167
1118	response. This is your gold standard.		1168
1119	3. **Compare:** Compare the "Generated	QUERY	1169
1120	Answer" against the "Ground Truth	Who was the second person to walk on the	1170
		moon?	1171

1172			1223
1173	GROUND TRUTH ANSWER	Here are the "Query", the "Ground Truth	1224
1174	Neil Armstrong was the first; Buzz Aldrin	Answer", and the "Generated Answer".	1225
1175	was second.		1226
1176		QUERY	1227
1177	GENERATED ANSWER	{query}	1228
1178	The second man on the moon was Michael		1229
1179	Collins.	GROUND TRUTH ANSWER	1230
1180		{ground_truth_answer}	1231
1181	YOUR RESPONSE		1232
1182	<think> Step 1: The query asks for the	GENERATED ANSWER	1233
1183	second person on the Moon. Step 2: The	{generated_answer}	1234
1184	ground truth identifies this person as		1235
1185	Buzz Aldrin. Step 3: The generated	YOUR RESPONSE	1236
1186	answer incorrectly identifies the		
1187	person as Michael Collins. Step 4:		
1188	This is a direct factual contradiction		
1189	with the ground truth. Therefore, the		
1190	answer is incorrect. </think>		
1191	<answer> incorrect </answer>		
1192			
1193	**Example 3**		
1194			
1195	QUERY		
1196	What are Newton's first two laws of motion?		
1197			
1198	GROUND TRUTH ANSWER		
1199	Newton's first law states that an object		
1200	will not change its motion unless a		
1201	force acts on it. The second law		
1202	states that the force on an object is		
1203	equal to its mass times its		
1204	acceleration.		
1205			
1206	GENERATED ANSWER		
1207	Newton's first law is the law of inertia,		
1208	stating an object in motion stays in		
1209	motion.		
1210			
1211	YOUR RESPONSE		
1212	<think> Step 1: The query asks for		
1213	Newton's first AND second laws. Step		
1214	2: The ground truth provides both		
1215	laws. Step 3: The generated answer		
1216	only provides the first law. Step 4:		
1217	The generated answer is incomplete as		
1218	it omits a critical part of the		
1219	information required by the query and		
1220	present in the ground truth.		
1221	Therefore, it is incorrect. </think>		
1222	<answer> incorrect </answer>		
		Appendix C. Prompts Used in the	1237
		Surgical Information	1238
		Assistant	1239
		This appendix contains all the prompts used in the Sur-	1240
		gical Information Assistant codebase. These prompts	1241
		are used to guide the language models in performing	1242
		various tasks such as question decomposition, informa-	1243
		tion retrieval, answer synthesis, searching a fallback	1244
		data-store, best-effort answer generation, and follow-	1245
		up question generation.	1246
		C.1. Question Decomposition Prompt	1247
		This prompt is used to break down complex surgical	1248
		questions into simpler sub-questions:	1249
		You are an expert at breaking complex	1250
		surgical questions into simpler ones.	1251
		Break the following question into	1252
		smaller sub-questions:	1253
			1254
		Question: {question}	1255
			1256
		Each sub-question should be independent	1257
		and answerable on it's own without	1258
		needing reference to other	1259
		sub-questions. Think of at least 3	1260
		sub-questions but no more than 7.	1261
			1262
		Think step-by-step and make sure to reason	1263
		through how break the question into	1264
		sub-questions.	1265
			1266

1267 Create new sub-questions in the following
 1268 format but do NOT answer the question.
 1269 Respond in the following format:
 1270
 1271 <think> Your reasoning here... </think>
 1272 <sub-question> The first sub-question...
 1273 </sub-question>
 1274 <sub-question> The second sub-question...
 1275 </sub-question>
 1276 ...
 1277 <sub-question> The last sub-question...
 1278 </sub-question>

1279 C.2. Answer Generation from Context 1280 Prompt

1281 This prompt is used to generate answers from retrieved
 1282 context for each sub-question. This results in multiple
 1283 LLM calls, but can be run asynchronously:

1284
 1285 Based on the given question and context,
 1286 generate an answer.
 1287 Question: {question}
 1288 Context: {context}

1289
 1290 Think step-by-step and make sure to reason
 1291 through how to generate an answer.
 1292 ONLY rely on the given context to
 1293 generate the answer.

1294
 1295 Include snippets of the context that
 1296 support your answer. Do NOT use any
 1297 information outside of the given
 1298 context to generate the answer.

1299
 1300 Respond in the following format:

1301
 1302 <think> Your reasoning here... </think>
 1303 <answer> The generated answer... </answer>
 1304 <snippet> First relevant snippet from the
 1305 context... </snippet>
 1306 <snippet> Second relevant snippet from the
 1307 context... </snippet>
 1308 ...
 1309 <snippet> The last relevant snippet from
 1310 the context </snippet>

C.3. Prompt for answer generation without context

This prompt is used to generate answers for sub-
 questions when no context is available (for ablation
 experiments).

You are a medical expert specializing in
 surgery. Answer the following question
 using your knowledge of surgical
 procedures, anatomy, and medical
 practices.

Question: {question}

Think step-by-step and provide a
 comprehensive answer based on your
 medical knowledge. If you're uncertain
 about any aspect, please indicate that
 in your response.

Respond in the following format:

<think> Your reasoning here... </think>
 <answer> The generated answer based on
 your medical knowledge... </answer>
 <confidence> High/Medium/Low - your
 confidence level in this answer
 </confidence>

C.4. Answer Synthesis Prompt

This prompt is used to synthesize answers from multi-
 ple sub-questions or route for further iteration in the
 DeRetSyn system.

You are a reasoning engine. Given the
 following sub-question answers,
 determine whether they are enough to
 fully answer the original question.
 ONLY rely on the knowledge to
 determine whether the question can be
 answered.

If yes, then provide the answer. Make your
 answer detailed and structured with
 sections, as appropriate. Include as
 much relevant information as possible
 from the knowledge provided.

If you determine that you cannot answer
 the original question, then suggest

1358	what additional questions should be		1404
1359	asked.	Question: {question}	1405
1360			1406
1361	Original Question:	Wikipedia Contexts:	1407
1362	{original_question}	{contexts}	1408
1363			1409
1364	Knowledge:	Think step-by-step to reason through your	1410
1365	{answers}	answer and consider the relevant	1411
1366		information from the contexts. Respond	1412
1367	Think step-by-step to reason through you	in the following format:	1413
1368	answer and consider the relevant	<think> Your reasoning here... </think>	1414
1369	information. Respond in the following	<answer> The synthesized answer...	1415
1370	format:	</answer>	1416
1371	<think> Your reasoning here... </think>		
1372	<can_answer> yes OR no </can_answer>		
1373	<answer> The answer to the original	C.7. Best Effort Answer Generation Prompt	1417
1374	question... </answer>	This prompt is used when the system needs to generate	1418
1375	<new_questions> The first new	a best-effort answer using Wikipedia:	1419
1376	sub-question... </new_questions>		
1377	<new_questions> The second new	You are a reasoning engine. Given the	1420
1378	sub-question... </new_questions>	following original question and	1421
1379	...	sub-question answers, formulate an	1422
1380	<new questions> The last new sub-question	answer to the best of your ability.	1423
1381	</new_questions>		1424
		Original Question:	1425
		{original_question}	1426
1382	C.5. Wikipedia Search Prompt		1427
1383	This prompt is used to generate search queries for	Knowledge:	1428
1384	Wikipedia:	{state["answers"]}	1429
		{state["wikipedia_results"]}	1430
1385	Given the following question, generate 3		1431
1386	search queries that would help find	Think step-by-step to reason through you	1432
1387	relevant information on Wikipedia. The	answer and consider the relevant	1433
1388	queries should be specific and focused	information. Respond in the following	1434
1389	on the key concepts in the question.	format:	1435
1390		<think> Your reasoning here... </think>	1436
1391	Question: {question}	<answer> The answer to the original	1437
1392		question... </answer>	1438
1393	Respond in the following format:		
1394	<query>first search query</query>	C.8. Follow-up Question Generation Prompt	1439
1395	<query>second search query</query>	This prompt is used to generate follow-up questions:	1440
1396	<query>third search query</query>		
		You are a reasoning engine. Given the	1441
1397	C.6. Wikipedia Context Synthesis Prompt	following original question and final	1442
1398	This prompt is used to synthesize information from	answer, generate 3 follow-up questions	1443
1399	Wikipedia contexts:	that help expand on the original	1444
1400	You are a reasoning engine. Given the	question and the answer in a step-wise	1445
1401	following question and Wikipedia	manner.	1446
1402	contexts, synthesize the information		1447
1403	to provide a comprehensive answer.	Original Question:	1448
		{original_question}	1449

1450		not, the LLM answers with general knowledge ac-	1494
1451	Final Answer:	knowledging it may hallucinate and directs the user	1495
1452	{final_answer}	to ask a more relevant question:	1496
1453			
1454	Think step-by-step to reason through your	Determine if the following question	1497
1455	answer and consider the relevant	requires access to specific medical	1498
1456	information. Respond in the following	documents related to detailed surgical	1499
1457	format:	information to be answered accurately.	1500
1458	<think> Your reasoning here... </think>	Think step-by-step and reason through your	1501
1459	<follow_up_questions> follow-up question	answer. Respond in the following	1502
1460	here... </follow_up_questions>	format:	1503
1461	<follow_up_questions> follow-up question		1504
1462	here... </follow_up_questions>	<thinking> Your reasoning here...	1505
1463	<follow_up_questions> follow-up question	</thinking>	1506
1464	here... </follow_up_questions>	<answer> yes OR no </answer>	1507
			1508
1465	C.9. Chain-of-Thought Generation Prompt	Here are some examples:	1509
1466	This prompt is used to generate detailed reasoning		1510
1467	for answers:	Question:	1511
1468	You are a reasoning engine. Based on the	What is the primary purpose of the	1512
1469	following question and knowledge,	coronary artery bypass graft?	1513
1470	provide a detailed, step-by-step	Response:	1514
1471	reasoning to arrive at an answer.	<thinking> The question is asking about a	1515
1472	Include at least 3 steps, but more as	coronary bypass graft which is related	1516
1473	needed.	to surgery. So yes, this question is	1517
1474		about surgery. </thinking>	1518
1475	Question:	<answer> yes </answer>	1519
1476	{state["original_question"]}		1520
1477		Question:	1521
1478	Knowledge:	Is machine learning useful for solving	1522
1479	{state["answers"]}	complex medical problems?	1523
1480	{state["wikipedia_results"]} if	Response:	1524
1481	"wikipedia_results" in state else ""}	<thinking> The question is asking about	1525
1482		machine learning as a tool for solving	1526
1483	Provide your response in this format:	complex medical problems. While the	1527
1484		question does specify medical	1528
1485	<think> Your reasoning here... </think>	problems, it does not refer to surgery	1529
1486	<answer> The final answer here... </answer>	or topics related to surgery.	1530
		</thinking>	1531
		<answer> no </answer>	1532
			1533
1487	Appendix D. Chat interaction	Question:	1534
1488	management prompts	Can you elaborate on your previous	1535
1489	The following prompts are used to manage the chat-	response about suturing or rephrase it?	1536
1490	interface with the Surgical Information Assistant.	Response:	1537
1491	D.1. Surgery Topic Classification Prompt	<thinking> The question is asking for an	1538
1492	This prompt is used to determine if a question is about	explanation of suturing that was given	1539
1493	surgery. If so, the DeRetSyn system is triggered. If	earlier in the conversation. While the	1540
		question does mention suturing which	1541
		is related to surgery, it is not	1542
		asking about suturing specifically but	1543

1544 rather requesting to explain a
1545 previous response. </thinking>
1546 <answer> no </answer>
1547
1548 Question:
1549 What are some core differences between
1550 robotic and laparoscopic inguinal
1551 hernia?
1552 Response:
1553 <thinking> The question is asking about
1554 the differences between robotic and
1555 laparoscopic inguinal hernia. Both
1556 robotic and laparoscopic surgery
1557 involve surgery on the inguinal
1558 hernia. This is clearly a question
1559 about surgery and requires
1560 investigation into documentation to
1561 answer. </thinking>
1562 <answer> yes </answer>
1563
1564 Question:
1565 {question}
1566 Response:

1567 **D.2. Context-Aware Input Update Prompt**

1568 This prompt is used to update user input with context
1569 from chat history:

1570 Given the following chat history and a new
1571 user input, update the user input to
1572 include any relevant context from the
1573 chat history. Only add necessary
1574 context, and keep the updated input
1575 concise.
1576 Think step-by-step and reason through your
1577 answer. Respond in the following
1578 format:
1579
1580 <thinking> Your reasoning here...
1581 </thinking>
1582 <updated_input> Updated user input here...
1583 </updated_input>
1584
1585 Here are some examples:

1586
1587 CHAT HISTORY:
1588 User: What are the main types of sutures
1589 used in surgery?
1590 Assistant: The main types of sutures used
1591 in surgery are absorbable and

non-absorbable sutures. Absorbable
sutures dissolve over time, while
non-absorbable sutures need to be
removed. Common absorbable sutures
include catgut and synthetic materials
like polyglycolic acid. Non-absorbable
sutures are often made of materials
like nylon or silk.

User: How long do absorbable sutures
typically last?
Assistant: The duration of absorbable
sutures varies depending on the
material used. Generally, they can
last anywhere from a few days to
several months. For example,
fast-absorbing sutures may dissolve
within 7-10 days, while
slower-absorbing sutures can take
60-90 days or even longer.

NEW USER INPUT:
What about non-absorbable ones?
RESPONSE:
<thinking> Considering the previous
discussion the user is referring to
non-absorbable sutures. Specifically,
they are asking about the lifespan or
removal time-frame for non-absorbable
sutures. I'll restate the question to
include this context from the chat
history. </thinking>
<updated_input> What is the typical
lifespan or removal timeframe for
non-absorbable sutures used in
surgery? </updated_input>

CHAT HISTORY:
User: What is a laparotomy?
Assistant: A laparotomy is a surgical
procedure that involves making a large
incision in the abdomen to gain access
to the abdominal cavity. It allows
surgeons to examine the abdominal
organs, diagnose problems, and perform
necessary treatments. Laparotomies are
used for various purposes, including
exploratory surgery, treatment of
abdominal trauma, and removal of
tumors.

1642	User: Are there any alternatives to this	if the surgical instrument is inserted	1692
1643	procedure?	directly into the abdominal cavity.	1693
1644	Assistant: Yes, there are alternatives to	3. Pain: Laparoscopic surgery can be	1694
1645	laparotomy, particularly minimally	painful, especially for patients with	1695
1646	invasive techniques. The main	pre-existing conditions or those who	1696
1647	alternative is laparoscopy, also known	have had previous laparoscopic	1697
1648	as keyhole surgery. In laparoscopy,	surgeries.	1698
1649	several small incisions are made		1699
1650	instead of one large incision. A	NEW USER INPUT:	1700
1651	camera and specialized instruments are	Can you think of any more?	1701
1652	inserted through these small incisions		1702
1653	to perform the surgery. This technique	RESPONSE:	1703
1654	often results in less pain, faster	<thinking> The user is asking if there are	1704
1655	recovery, and smaller scars compared	any more complications, but it's	1705
1656	to traditional laparotomy.	important to note that the original	1706
1657		question was about laparoscopic	1707
1658	NEW USER INPUT:	surgery, not about potential	1708
1659	What are the risks?	complications. I'll restate the	1709
1660		question to clarify that the original	1710
1661	RESPONSE:	topic was laparoscopic surgery.	1711
1662	<thinking> The user is asking about risks,	</thinking>	1712
1663	but it's not clear whether they're	<updated_input> Are there any more	1713
1664	referring to laparotomy or	potential complications associated	1714
1665	laparoscopy, both of which were	with laparoscopic surgery besides	1715
1666	discussed in the previous messages.	infection, stool bleeding, and pain?	1716
1667	Since laparotomy was the original	</updated_input>	1717
1668	topic and laparoscopy was introduced		1718
1669	as an alternative, it would be helpful		1719
1670	to ask about the risks of both		1720
1671	procedures for a comprehensive answer.	CHAT HISTORY:	1721
1672	</thinking>	{formatted_history}	1722
1673	<updated_input> What are the risks		1723
1674	associated with both laparotomy and	New User Input:	1724
1675	laparoscopy procedures?	{user_input}	1725
1676	</updated_input>		1726
1677		RESPONSE:	1727
1678			
1679			
1680	CHAT HISTORY:	Appendix E. Prompts used to test	1728
1681	User: What are the potential complications	alternative RAG	1729
1682	of laparoscopic surgery?	paradigms	1730
1683	Assistant: Laparoscopic surgery can have	E.1. Chain-of-thought prompting without	1731
1684	several potential complications,	context	1732
1685	including:		
1686	1. Infection: Laparoscopic surgery can	You are a medical expert. Please answer	1733
1687	lead to infections, especially in	the following question:	1734
1688	patients with infections that can be	Think step-by-step and provide a detailed	1735
1689	spread through the abdominal cavity.	reasoning process to arrive at your	1736
1690	2. Stool bleeding: Laparoscopic surgery	answer. Include at least 3 steps in	1737
1691	can lead to stool bleeding, especially	your reasoning, but more as needed.	1738

1739		3. Provide comprehensive and accurate	1785
1740	Respond in the following format:	answers based on the retrieved content	1786
1741		4. If the information is not available in	1787
1742	<think> Your reasoning here... </think>	the search results, acknowledge the	1788
1743	<answer> Your final answer here...	limitations	1789
1744	</answer>	5. Think step-by-step to reason through	1790
1745		complex questions	1791
1746	Question: {question}	6. Cite specific parts of the retrieved	1792
		documents when appropriate	1793
1747	E.2. Chain-of-thought prompting with	7. Focus on providing factual medical	1794
1748	context	information rather than opinions	1795
			1796
1749	You are a medical assistant specializing	Remember accuracy is crucial. Provide all	1797
1750	in surgical information. Use the	reasoning and the final answer.	1798
1751	following context to answer the		
1752	question.		
1753	If you cannot find the answer in the	Appendix F. Prompt Design	1799
1754	context, say "I don't have enough	Principles	1800
1755	information to answer this question."		
1756		The prompts used in the Surgical Information Assis-	1801
1757	Context:	tant follow several key design principles:	1802
1758	{context_docs}		
1759			
1760	Question: {question}	F.1. Step-by-Step Reasoning	1803
1761			
1762	Think step-by-step to reason through your	All prompts encourage the language model to think	1804
1763	answer and consider the relevant	step-by-step and provide detailed reasoning before	1805
1764	information from the context. Respond	arriving at an answer. This approach, often referred	1806
1765	in the following format:	to as chain-of-thought prompting, has been shown	1807
1766	<think> Your reasoning here... </think>	to improve the accuracy and reliability of language	1808
1767	<answer> The answer to the question...	model outputs, especially for complex tasks.	1809
1768	</answer>		
		F.2. Structured Output Format	1810
1769	E.3. ReAct Agent System Prompt		
1770	This prompt is used for the ReAct agent in the evalu-	The prompts use a consistent structured output for-	1811
1771	ation scripts:	mat with XML-like tags (e.g., <think>, <answer>,	1812
1772	You are a medical assistant specializing	<sub-question>) to clearly separate different compo-	1813
1773	in surgical information. Your goal is	nents of the response. This structured format makes	1814
1774	to provide accurate and helpful	it easier to parse and extract specific information from	1815
1775	information about surgical procedures,	the language model's output.	1816
1776	techniques, and related medical		
1777	knowledge.	F.3. Few-Shot Examples	1817
1778			
1779	When answering questions:	Several prompts include few-shot examples to demon-	1818
1780	1. Use the DocumentSearch tool to find	strate the expected reasoning process and output for-	1819
1781	relevant information in the surgical	mat. These examples help guide the language model	1820
1782	database	to produce responses that follow the desired pattern	1821
1783	2. Analyze the retrieved information	and meet the quality standards.	1822
1784	carefully	F.4. Task-Specific Instructions	1823
		Each prompt contains clear, task-specific instructions	1824
		that define the role of the language model (e.g., "You	1825

are an expert at breaking complex surgical questions into simpler ones”) and the specific task it needs to perform. This helps focus the model’s attention on the relevant aspects of the task.

F.5. Context Awareness

The prompts are designed to make the language model aware of the context in which it operates, including the original question, retrieved information, and previous interactions. This context awareness is crucial for generating coherent and relevant responses.

Appendix G. Example responses from different experimental designs

G.1. Question 1

Why should soiled instruments be cleaned within twenty minutes of contamination?

Known answer: Soiled instruments should be cleaned within twenty minutes because stainless steel can experience corrosion from biological and chemical liquids.

Llama32 Deretsyn: The soiled instruments should be cleaned within twenty minutes of contamination to prevent the potential risks associated with them, including corrosion, injury from sharp instruments, and the spread of infection due to the presence of biological and chemical liquids. This time frame is significant because it is the time frame within which organic matter such as blood or tissue can be prevented from drying on the surface of the instrument. If not cleaned within this time frame, the consequences can include damage to the instruments, corrosion, and the potential for injury. Furthermore, there are specific guidelines that dictate the cleaning of soiled instruments within twenty minutes of contamination.

Llama32 RAG: Soiled instruments should be cleaned within twenty minutes of contamination to prevent the growth of microorganisms and the drying of

organic matter, which can lead to the formation of biofilms and make the instrument more difficult to clean and disinfect.

Llama32 React: Soiled instruments should be cleaned within twenty minutes of contamination to prevent the growth and spread of microorganisms. This is because microorganisms can multiply rapidly, especially in warm, moist environments, and can cause infections if not removed promptly.

GPT40: Soiled instruments should be cleaned within twenty minutes of contamination to prevent organic material from drying and adhering to the surface, to reduce the risk of biofilm formation that can make sterilization less effective, to minimize the risk of cross-contamination and infection, to maintain the integrity and functionality of the instruments, and to adhere to best practices in infection control.

Llama31: Soiled instruments should be cleaned within twenty minutes of contamination to prevent bacterial growth, biofilm formation, and subsequent infection transmission. This timeframe allows for effective cleaning and disinfection, reducing the risk of hospital-acquired infections and ensuring patient safety.

G.2. Question 2:

What factors can put the liver remnant at risk for further ischemic injury postoperatively after liver resection surgery?

Known answer: Hypoxia, hypotension, septicemia, or drug toxicity can put the liver remnant at risk for further ischemic injury postoperatively.

Llama32 Deretsyn: The factors that can put the liver remnant at risk for further ischemic injury postoperatively include:

1. **Hemodynamic instability**: This can affect blood flow to the liver remnant

1921	and increase the risk of ischemic	liver remnant and increase the risk of	1972
1922	injury.	ischemic injury.	1973
1923	2. **Hyperdynamic circulation** : This can	15. **Postoperative factors** : Hypoxia,	1974
1924	also affect blood flow to the liver	hypotension, septicemia, drug	1975
1925	remnant and increase the risk of	toxicity, hyperdynamic circulation,	1976
1926	ischemic injury.	and hyperglycemia can all contribute	1977
1927	3. **Septicemia** : This is a potential	to ischemic injury in the liver	1978
1928	cause of ischemic injury in the liver	remnant postoperatively.	1979
1929	remnant postoperatively.	16. **Surgical procedure** : The surgical	1980
1930	4. **Drug toxicity** : This is another	procedure itself can increase the risk	1981
1931	potential cause of ischemic injury in	of hypoxia and hypotension, prolong	1982
1932	the liver remnant postoperatively.	warm ischemia time, and cause	1983
1933	5. **Volume expansion** : This can affect	hemodynamic instability, all of which	1984
1934	blood flow to the liver remnant and	can increase the risk of ischemic	1985
1935	increase the risk of ischemic injury.	injury in the liver remnant.	1986
1936	6. **Glycemic control** : Poor glycemic		1987
1937	control can lead to hyperglycemia,	These factors can all contribute to	1988
1938	which can contribute to ischemic	ischemic injury in the liver remnant	1989
1939	injury in the liver remnant.	postoperatively, either by affecting	1990
1940	7. **Liver regeneration** : While liver	blood flow to the liver remnant or by	1991
1941	regeneration is a natural process, it	increasing the risk of hypoxia,	1992
1942	can also lead to changes in blood flow	hypotension, septicemia, and drug	1993
1943	to the liver remnant, increasing the	toxicity.	1994
1944	risk of ischemic injury.	Llama32 RAG: Hypoxia, hypotension,	1995
1945	8. **Ascites** : This can put pressure on	septicemia or drug toxicity, and	1996
1946	the liver remnant and impair venous	hyperdynamic circulation due to	1997
1947	return to the heart, increasing the	cirrhosis or impaired liver function	1998
1948	risk of ischemic injury.	can put the liver remnant at risk for	1999
1949	9. **Bleeding** : This can lead to	further ischemic injury	2000
1950	hemodynamic instability and affect	postoperatively.	2001
1951	blood flow to the liver remnant,	Llama32 React: The following factors can	2002
1952	increasing the risk of ischemic injury.	put the liver remnant at risk for	2003
1953	10. **Portal hypertension** : This can	further ischemic injury	2004
1954	affect blood flow to the liver remnant	postoperatively: 1. Warm ischemia	2005
1955	and increase the risk of ischemic	time, 2. Fatty or cirrhotic livers, 3.	2006
1956	injury.	Large resections, 4. Prolonged	2007
1957	11. **Cardiac function** : Poor cardiac	surgical time, 5. Blood loss and	2008
1958	function can lead to hypotension and	cardiovascular instability, 6.	2009
1959	affect blood flow to the liver	Non-steroidal anti-inflammatory drugs	2010
1960	remnant, increasing the risk of	(NSAIDs), 7. Paracetamol	2011
1961	ischemic injury.	(acetaminophen) use, 8. Hypoglycemia,	2012
1962	12. **Coagulation status** : Abnormal	9. Hypothermia, 10. Poor patient	2013
1963	coagulation status can lead to	selection.	2014
1964	bleeding and affect blood flow to the	GPT40: Factors that can put the liver	2015
1965	liver remnant, increasing the risk of	remnant at risk for further ischemic	2016
1966	ischemic injury.	injury postoperatively include	2017
1967	13. **Warm ischemia time** : Prolonged warm	vascular and perfusion issues (e.g.,	2018
1968	ischemia time can increase the risk of	portal vein thrombosis, hepatic artery	2019
1969	ischemic injury in the liver remnant.	injury, small-for-size syndrome),	2020
1970	14. **Liver function** : Poor liver	systemic and local hypoxia,	2021
1971	function can affect blood flow to the	inflammatory and oxidative stress	2022

related to ischemia-reperfusion injury, pre-existing conditions that impair liver function (e.g., cirrhosis, steatosis, diabetes), poor postoperative hemodynamic management (e.g., hypotension, inadequate resuscitation), and a small or insufficiently functional liver remnant that cannot meet metabolic demands.

Llama31: The liver remnant is at risk for further ischemic injury postoperatively due to a combination of factors, including:

- Hemodynamic factors (e.g., changes in blood flow, blood pressure, or cardiac output)
- Surgical factors (e.g., extent of liver resection, presence of major blood vessels in the area of resection)
- Patient-related factors (e.g., overall health, age, presence of comorbidities)
- Specific conditions (e.g., small liver remnant, portal hypertension, biliary obstruction, dehydration or hypovolemia)
- Inadequate intraoperative and postoperative care (e.g., fluid management, blood pressure management, monitoring and management of complications)